

CLAIMS

What is claimed is:

1 1. A method for performing a distributed transaction in a shared-nothing database
2 system, the method comprising:
3 on a first shared-nothing node of said shared-nothing database system, causing a
4 coordinator that is coordinating the distributed transaction to store information
5 that indicates status of said distributed transaction on a persistent storage
6 device;
7 wherein the persistent storage device is accessible to a participant that is to perform
8 one or more operations as part of said distributed transaction;
9 wherein the participant resides on a second shared-nothing node of said shared-
10 nothing database system; and
11 on the second shared-nothing node of said shared-nothing database system, causing
12 the participant to determine the status of said distributed transaction by
13 reading the information from the persistent storage device.

1 2. The method of Claim 1 wherein:
2 the participant is a first participant of a plurality of participants in said distributed
3 transaction;
4 the plurality of participants includes a second participant that does not have access to
5 said persistent storage device; and
6 the method further comprises the step of the coordinator interacting with the second
7 participant according to a two-phase commit protocol.

1 3. The method of Claim 1 further comprising the steps of:

2 the coordinator committing the distributed transaction;
3 after the coordinator commits the distributed transaction, the coordinator sending a
4 commit message to the participant; and
5 preventing the information that indicates the status of the distributed transaction from
6 being overwritten or deleted until a set of conditions is satisfied, wherein one
7 condition in said set of conditions is that the coordinator receives a commit
8 acknowledge message from said participant.

1 4. The method of Claim 1 further comprising the steps of:
2 the participant sending a first piece of information to the coordinator, wherein the first
3 piece of information is associated with work performed by said participant as
4 part of said distributed transaction; and
5 the coordinator performing a comparison between the first piece of information and
6 information associated with a redo log of said second shared-nothing node;
7 and
8 the coordinator determining whether to commit the transaction based, at least in part,
9 on said comparison.

1 5. The method of Claim 4 wherein piece of information includes a log-sequence-number
2 of the latest change made by the participant as part of the distributed transaction.

1 6. The method of Claim 5 wherein the step of sending includes the steps of:
2 the participant identifying a message that is being sent to said first shared-nothing
3 node for a purpose unrelated to the distributed transaction; and
4 piggybacking the log-sequence number on said message.

1 7. A method for performing a distributed transaction in a shared-nothing database
2 system, the method comprising:
3 assigning a participant to perform one or more operations as part of said distributed
4 transaction;
5 wherein the participant resides on a first shared-nothing node of said shared-nothing
6 system;
7 causing said participant to store, on a persistent storage device, status information that
8 indicates changes made by the participant during performance of said one or
9 more operations;
10 wherein the persistent storage device is accessible to a coordinator that is responsible
11 for coordinating said distributed transaction;
12 wherein the coordinator resides on a second shared-nothing node of said shared-
13 nothing database system;
14 on said second shared-nothing node of said shared-nothing database system, causing
15 said coordinator to determine, based on the status information on said
16 persistent storage device, whether the participant has written to persistent
17 storage changes produced by performance of the one or more operations; and
18 the coordinator process determining whether the distributed transaction can be
19 committed based, at least in part, on whether the participant has written to
20 persistent storage changes produced by performance of the one or more
21 operations.

1 8. The method of Claim 7 wherein:

2 the step of causing said participant to store, on a persistent storage device, status
3 information that indicates changes made by the participant during
4 performance of said one or more operations includes
5 causing said participant to store redo information in a redo log on said
6 persistent storage device; and

7 the step of causing said coordinator to determine, based on the status information on
8 said persistent storage device, whether the participant has written to persistent
9 storage changes produced by performance of the one or more operations
10 includes
11 inspecting the redo log of the participant to determine whether the redo
12 information for said changes have been written to said persistent
13 storage.

1 9. The method of Claim 7 wherein:

2 the participant is a first participant of a plurality of participants in said distributed
3 transaction;

4 the plurality of participants includes a second participant that stores status
5 information on a second persistent storage device that is not accessible by said
6 coordinator; and

7 the method further comprises the step of the coordinator interacting with the second
8 participant according to a two-phase commit protocol.

1 10. The method of Claim 7 wherein:

2 the information on said persistent storage device indicates that the participant has not
3 written to persistent storage changes produced by performance of the one or
4 more operations; and
5 the method further comprises the coordinator sending a force redo message to the
6 participant to cause the participant to write to persistent storage the changes
7 produced by performance of the one or more operations.

1 11. The method of Claim 10 wherein the step of sending a force redo message includes
2 the steps of:
3 identifying a message that is being sent to said first shared-nothing node for a purpose
4 unrelated to the distributed transaction; and
5 piggybacking the force redo message on said message.

1 12. A computer-readable medium carrying one or more sequences of instructions which,
2 when executed by one or more processors, causes the one or more processors to perform the
3 method recited in Claim 1.

1 13. A computer-readable medium carrying one or more sequences of instructions which,
2 when executed by one or more processors, causes the one or more processors to perform the
3 method recited in Claim 2.

1 14. A computer-readable medium carrying one or more sequences of instructions which,
2 when executed by one or more processors, causes the one or more processors to perform the
3 method recited in Claim 3.

1 15. A computer-readable medium carrying one or more sequences of instructions which,
2 when executed by one or more processors, causes the one or more processors to perform the
3 method recited in Claim 4.

1 16. A computer-readable medium carrying one or more sequences of instructions which,
2 when executed by one or more processors, causes the one or more processors to perform the
3 method recited in Claim 5.

1 17. A computer-readable medium carrying one or more sequences of instructions which,
2 when executed by one or more processors, causes the one or more processors to perform the
3 method recited in Claim 6.

1 18. A computer-readable medium carrying one or more sequences of instructions which,
2 when executed by one or more processors, causes the one or more processors to perform the
3 method recited in Claim 7.

1 19. A computer-readable medium carrying one or more sequences of instructions which,
2 when executed by one or more processors, causes the one or more processors to perform the
3 method recited in Claim 8.

1 20. A computer-readable medium carrying one or more sequences of instructions which,
2 when executed by one or more processors, causes the one or more processors to perform the
3 method recited in Claim 9.

1 21. A computer-readable medium carrying one or more sequences of instructions which,
2 when executed by one or more processors, causes the one or more processors to perform the
3 method recited in Claim 10.

1 22. A computer-readable medium carrying one or more sequences of instructions which,
2 when executed by one or more processors, causes the one or more processors to perform the
3 method recited in Claim 11.